

2

BACKGROUND AND STATE OF THE ART

You must have a map, no matter how rough. Otherwise you wander all over the place.

— J.R.R. Tolkien

THIS chapter discusses the state-of-the-art in the areas of WebAssembly and Software Diversification. In Section 2.1 we discuss WebAssembly, focusing on its design and security model. Besides, we discuss the current state-of-the-art of WebAssembly research. In Section 2.2 we discuss related works in the area of Software Diversification. Moreover, we delve into the open challenges regarding the diversification of WebAssembly programs.

2.1 WebAssembly

The W3C publicly announced the WebAssembly (Wasm) language in 2017 as the fourth scripting language supported in all major web browser vendors. WebAssembly is a binary instruction format for a stack-based virtual machine and was officially consolidated by the work of Haas et al. [7] in 2017 and extended by Rossberg et al. in 2018 [38]. It is designed to be fast, portable, self-contained, and secure.

Moreover, WebAssembly has been evolving outside web browsers since its first announcement. Previous works demonstrated that using WebAssembly as an intermediate layer is better in terms of startup time and memory usage than containerization and virtualization [10, 11]. Consequently, in 2019, the Bytecode Alliance proposed WebAssembly System Interface (WASI) [39]. WASI pioneered the execution of WebAssembly with a POSIX system interface protocol, making it possible to execute Wasm closer to the underlying operating system. Therefore, it standardizes the adoption of WebAssembly in heterogeneous platforms [40], i.e., IoT and Edge computing [41, 42].

⁰Compilation probe time 2023/11/23 13:19:08

Currently, WebAssembly serves a variety of functions in browsers, ranging from gaming to cryptomining [43]. Other applications include text processing, visualization, media processing, programming language testing, online gambling, bar code and QR code fast reading, hashing, and PDF viewing. On the backend, WebAssembly notably excels in short-running tasks. As such, it is particularly suitable for Function as a Service (FaaS) platforms like Cloudflare and Fastly. The subsequent text in this chapter focuses specifically on WebAssembly version 1.0. However, the tools, techniques, and methodologies discussed are applicable to future WebAssembly versions.

2.1.1 From source code to WebAssembly

WebAssembly programs are compiled from source languages like C/C++, Rust, or Go ahead of time, which means that Wasm binaries can benefit from the optimizations of the source language compiler. The resulting WebAssembly program is like a traditional shared library, containing instruction codes, symbols, and exported functions. A host environment is in charge of complementing the Wasm program, such as providing external functions required for execution within the host engine. For instance, functions for interacting with an HTML page’s DOM are imported into the Wasm binary when invoked from JavaScript code in the browser.

In Listing 2.1 and Listing 2.2, we present a Rust program alongside its corresponding WebAssembly binary. The Rust program in Listing 2.1 iteratively calculates the Fibonacci sequence up to a given number that comes from the host engine. The code in the program encompasses various elements such as vector allocations, external function usage, and a function definition that includes a loop, conditionals, function calls, and memory accesses. The Wasm code shown in Listing 2.2 is simplified in its textual format, known as WAT¹. The function prototype in lines 4 and 5 of Listing 2.1 are converted into imported function, as seen in lines 8 and 9 of Listing 2.2. The `fibonacci` function, spanning lines 7 to 20 in Listing 2.1, is compiled into a Wasm function covering lines 14 to 31 in Listing 2.2. Within this function, the translation of various Rust language constructs into Wasm can be observed. For instance, the `for` loop found in line 14 of Listing 2.1 is mapped to a block structure in lines 17 to 31 of Listing 2.2. The breaking condition of the loop is transformed into a conditional branch, as depicted in line 23 of Listing 2.2. In this scenario, the function yields the final set value in the `local` variable. Note that for optimization purposes, the loop concludes by returning the result value, instead of returning post completion of the loop.

¹The WAT text format is primarily designed for human readability and for low-level manual editing.

```

1 ...
2 // Imported from host
3 extern "C" {
4     fn log(s: &str);
5     fn get_input() -> usize; }
6
7 fn fibonacci(n: usize) -> i32 {
8     // Iterative fibonacci
9     // Create a vector of size n+1
10    let mut fibo_result = vec![0; n + 1];
11    // Set ith 0 and 1
12    fibo_result[0] = 1;
13    fibo_result[1] = 1;
14    for i in 2..=n {
15        // f[i] = f[i-1] + f[i-2]
16        fibo_result[i] = fibo_result[i - 1] + fibo_result[i - 2];
17    }
18    // Return the last element
19    return fibo_result[n];
20 }
21 // Pub to export the function
22 pub fn main() {
23     // Get the input from the user
24     let nth = get_input();
25     // Calculate the fibonacci
26     let fib = fibonacci(get_input());
27     // Print the result in the host imported function
28     log(&format!("{}", fib));
29 }
```

Listing 2.1: Example Rust program which includes, external function usage, a function definition featuring a loop, function calls, imported functions, and memory accesses.

There exist several compilers that turn source code into WebAssembly binaries. For example, LLVM compiles to WebAssembly as a backend option since its 7.1.0 release in early 2019², supporting a diverse set of frontend languages like C/C++, Rust, Go, and AssemblyScript³. Significantly, a study by Hilbig [43] reveals that 70% of WebAssembly binaries are generated using LLVM-based compilers. The main advantage of using LLVM is that it provides a modular and state-of-the-art optimization infrastructure for WebAssembly binaries. Recently, the Kotlin Multiplatform framework ⁴ has incorporated WebAssembly as a compilation target, enabling the compilation of Kotlin code to WebAssembly.

²<https://github.com/llvm/llvm-project/releases/tag/llvmorg-7.1.0>

³A subset of the TypeScript language

⁴<https://kotlinlang.org/docs/wasm-overview.html>

```

1 ; WebAssembly magic bytes(\0asm) and version (1.0) ;
2 (module
3 ...
4 ; Type section: 0x01 0x00 0x00 0x00 0x13 ...
5 (type (;type index 0;) (func (param i32 i32)))
6 ...
7 ; Import section: 0x02 0x00 0x00 0x00 0x57 ...
8 (import "__wbg__" "__wbg_log" (func (;1;) (type 0)))
9 (import "__wbg__" "__wbg_getinput" (func (;2;) (type 8)))
10 ...
11 ; Custom section: 0x00 0x00 0x00 0x00 0x7E ;
12 (@custom "name" "...")
13 ...
14 (func (;func index 40;) (type 1) (param i32) (result i32)
15   (local i32 i32 i32 i32 i32) ;local variables;
16 ...
17   loop ; label = @1 ;
18 ...
19   i32.eqz
20   if ; label = @2 Compare the top of the stack ;
21 ...
22   local.get 0
23   return ; Return the last element which is saved in local 0 ;
24 end
25 ...
26 block ;label = @2 ;
27 ...
28   i32.store ; Store the fib value in the mem assigned to the
29   ↪ result array;
30   br 1 (;@1;) ;Continue the loop;
31 end
32 ...
33 (func (;44;) (type 8) (result i32)
34 ...
35   call 2 ; Calling the imported function to get input ;
36   i32.store ; Store the input in memory ;
37 ...
38 (func (;45;) (type 7)
39   (local i32 i32 i32)
40 ...
41   call 44
42   call 40 ; Calling fibo function ;
43   i32.store offset=20
44 ...
45 (table (;0;) 33 33 funcref)
46 ; Memory section: 0x05 0x00 0x00 0x00 0x03 ...
47 (memory (;0;) 17)
48 ; Global section: 0x06 0x00 0x00 0x00 0x11...
49 (global (;global index 0;) (mut i32 ;mut global;) (i32.const 1048576))
50 ...
51 ; Export section: 0x07 0x00 0x00 0x00 0x72 ...
52 (export "memory" (memory 0))
53 (export "fibo" (func 40))
54 (export "main" (func 45))
55 ...
56 ; Data section: 0x0d 0x00 0x00 0x03 0xEF ...
57 (data (;data segment index 0;) (i32.const 1048576) "invalid args...")
58 ...
59 ; Custom section: 0x00 0x00 0x00 0x00 0x2F ;
60 (@custom "producers" "...")
```

Listing 2.2: Refer to Listing 2.1 for the Rust code example. This example showcases the transition from Rust to Wasm, where numerous high-level language attributes convert into multiple Wasm instructions. For clarity, we've marked elements and portions of the WebAssembly binary as comments.

A recent trend in the WebAssembly ecosystem involves porting various programming languages by converting both the language's engine or interpreter and the source code into a WebAssembly program. For example, Javy⁵ encapsulates JavaScript code within the QuickJS interpreter, demonstrating that direct source code conversion to WebAssembly isn't always required. If an interpreter for a specific language can be compiled to WebAssembly, it allows for the bundling of both the interpreter and the language into a single, isolated WebAssembly binary. Similarly, Blazor⁶ facilitates the execution of .NET Common Intermediate Language (CIL) in WebAssembly binaries for browser-based applications. However, packaging the interpreter and the code in one single standalone WebAssembly binary is still immature and faces challenges. For example, the absence of JIT compilation for the "interpreted" code makes it less suitable for long-running tasks [44]. On the other hand, it proves effective for short-running tasks, particularly those executed in Edge-Cloud computing platforms.

2.1.2 WebAssembly's binary format

The Wasm binary format is close to machine code and already optimized, being a consecutive collection of sections. In Figure 2.1 we show the binary format of a Wasm section. A Wasm section starts with a 1-byte section ID, followed by a 4-byte section size, and concludes with the section content, which precisely matches the size indicated earlier. A WebAssembly binary contains sections of 13 types, each with a specific semantic role and placement within the module. For instance, the *Custom Section* stores metadata like the compiler used to generate the binary, while the *Type Section* contains function signatures that serve to validate the *Function Section*. The *Import Section* lists elements imported from the host, and the *Function Section* details the functions defined within the binary. Other sections like *Table*, *Memory*, and *Global Sections* specify the structure for indirect calls, unmanaged linear memories, and global variables, respectively. *Export*, *Start*, *Element*, *Code*, *Data*, and *Data Count Sections* handle aspects ranging from declaring elements for host engine access to initializing program state, declaring bytecode instructions per function, and initializing linear memory. Each of these sections must occur only once in a binary and can be empty. For clarity, we also annotate sections as comments in the Wasm code in Listing 2.2.

A WebAssembly binary can be processed efficiently due to its organization into a contiguous array of sections. For instance, this structure permits compilers to boost the compilation process either through parallel parsing. Moreover, the *Code Section*'s instructions are further compacted through the use of the LEB128⁷

⁵<https://github.com/bytocodealliance/javy>

⁶<https://dotnet.microsoft.com/en-us/apps/aspnet/web-apps/blazor>

⁷<https://en.wikipedia.org/wiki/LEB128>



Figure 2.1: Memory byte representation of a WebAssembly binary section, starting with a 1-byte section ID, followed by an 8-byte section size, and finally the section content.

encoding. Consequently, Wasm binaries are not only fast to validate and compile, but also swift to transmit over a network.

2.1.3 WebAssembly’s runtime

The WebAssembly’s runtime characterizes the behavior of WebAssembly programs during execution. This section describes the main components of the WebAssembly runtime, namely the execution stack, functions, memory model, and execution process. These components are crucial for understanding both the WebAssembly’s control-flow and the analysis of WebAssembly binaries.

Execution Stack: At runtime, WebAssembly engines instantiate a WebAssembly module. This module is a runtime representation of a loaded and initialized WebAssembly binary described in Section 2.1.2. The primary component of a module instance is its Execution Stack. The Execution Stack stores typed values, labels, and control frames. Labels manage block instruction starts and loop starts. Control frames manage function calls and function returns. Values within the stack can only be static types. These types include `i32` for 32-bit signed integers, `i64` for 64-bit signed integers, `f32` for 32-bit floats, and `f64` for 64-bit floats. Abstract types such as classes, objects, and arrays are not supported natively. Instead, these types are abstracted into primitive types during compilation and stored in linear memory.

Functions: At runtime, WebAssembly functions are closures over the module instance, grouping locals and function bodies. Locals are typed variables that are local to a specific function invocation. A function body is a sequence of instructions that are executed when the function is called. Each instruction either reads from the execution stack, writes to the execution stack, or modifies the control-flow of the function. Recalling the example WebAssembly binary previously showed, the local variable declarations and typed instructions that are evaluated using the stack can be appreciated between Line 12 and Line 38 in Listing 2.2. Each instruction reads its operands from the stack and pushes back the result. Notice that, numeric instructions are annotated with its corresponding type. In the case of Listing 2.2, the result value of the main function is the calculation of the last instruction, `i32.add` in line 38. As the listing also shows, instructions are annotated with a numeric type.

Memory model: A WebAssembly module instance incorporates three key types of memory-related components: linear memory, local variables and global variables. These components can either be managed solely by the host engine or shared with the WebAssembly binary itself. This division of responsibility is often categorized as *managed* and *unmanaged* memory [17]. Managed refers to components that are exclusively modified by the host engine at the lowest level, e.g. when the WebAssembly binary is JITed, while unmanaged components can also be altered through WebAssembly opcodes. First, modules may include multiple linear memory instances, which are contiguous arrays of bytes. These are accessed using 32-bit integers (`i32`) and are shareable only between the initiating engine and the WebAssembly binary. Generally, these linear memories are considered to be unmanaged, e.g., line 28 of Listing 2.2 shows an explicit memory access opcode. Second, there are global instances, which are variables accompanied by values and mutability flags (see example in line 49 of Listing 2.2). These globals are managed by the host engine, which controls their allocation and memory placement completely oblivious to the WebAssembly binary scope. They can only be accessed via their declaration index, prohibiting dynamic addressing. Third, local variables are mutable and specific to a given function instance (e.g., line 15 in Listing 2.2). They are accessible only through their index relative to the executing function and are part of the data managed by the host engine.

WebAssembly module execution: While a WebAssembly binary could be interpreted, the most practical approach is to JIT compile it into machine code [45]. The main reason is that WebAssembly is optimized and closely aligned with machine code, leading to swift JIT compilation for execution. Browser engines such as V8⁸ and SpiderMonkey⁹ utilize this strategy when executing WebAssembly binaries in browser clients¹⁰. Once JITed, the WebAssembly binary operates within a sandboxed environment, accessing the host environment exclusively through imported functions. This sandboxing follows the Software Fault Isolation(SFI) strategy, meaning that a WebAssembly program cannot arbitrarily access code or data of its runtime.

WebAssembly standalone engines: While initially intended for browsers, WebAssembly has undergone significant evolution, primarily due to WASI[39]. WASI establishes a standardized POSIX-like interface for interactions between WebAssembly modules and host environments. Compilers can generate WebAssembly binaries that implement WASI, which allows execution in standalone engines. These binaries can then be executed by standalone engines across a variety of environments, including the cloud, servers, and IoT devices

⁸<https://chromium.googlesource.com/v8/v8.git>

⁹<https://spidermonkey.dev/>

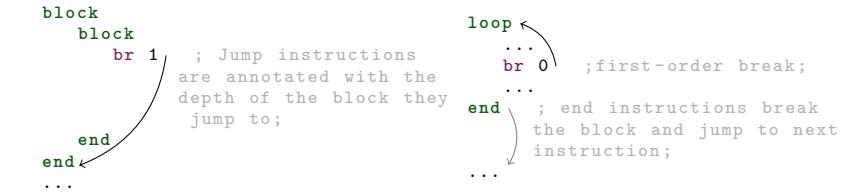
¹⁰In practice, browsers initially employ a baseline compiler to ensure the rapid availability of incoming source code. Simultaneously, the optimizing compiler operates in the background. Consequently, the first generated machine code is eventually supplanted by the optimized version. Given that the objectives align, we consider this two-fold process as one.

[46, 45]. Similarly to browsers, these engines often translate WebAssembly into machine code via JIT compilation, ensuring a sandboxed execution process. Standalone engines such as WASM¹¹, Wasmer¹², Wasmtime¹³, WAVM¹⁴, and Sledge[47] have been developed to support both WebAssembly and WASI.

2.1.4 WebAssembly’s control-flow

A WebAssembly function groups instructions into blocks, with the function’s entrypoint acting as the root block. In contrast to conventional assembly code, control-flow structures in Wasm leap between block boundaries rather than arbitrary positions within the code, effectively prohibiting `gotos` to random code positions. Each block may specify the needed execution stack state before execution as well as the resultant execution stack state once its instructions have been executed. Typically, the execution stack state is simply the quantity and numeric type of values on the stack. This stack state is used to validate the binary during compilation and to ensure that the stack is in a valid state before the execution of the block’s instructions. Blocks in Wasm are explicit (see instructions `block` and `end` in lines 16 and 34 of Listing 2.2), delineating where they start and end. By design, a block cannot reference or execute code from external blocks.

During runtime, WebAssembly break instructions can only jump to one of its enclosing blocks. Breaks, except for those within loop constructions, jump to the block’s end and continue to the next immediate instruction. For instance, after line 31 of Listing 2.2, the execution would proceed to line 32. Within a loop, the end of a block results in a jump to the block’s beginning, thus restarting the loop. For example, if line 29 of Listing 2.2 evaluates as false, the next instruction to be executed in the loop would be line 18. Listing 2.3 provides an example for better understanding, comparing a standard block and a loop block in a Wasm function.



Listing 2.3: Example of breaking a block and a loop in WebAssembly.

Each break instruction includes the depth of the enclosing block as an operand. This depth is used to identify the target block for the break instruction.

¹¹<https://github.com/wasm3/wasm3>

¹²<https://wasmer.io/>

¹³<https://github.com/bytocodealliance/wasmtime>

¹⁴<https://github.com/WAVM/WAVM>

For example, in the left-most part of the previously discussed listing, a break instruction with a depth of 1 would jump past two enclosing blocks. This design hardens the rewriting of WebAssembly binaries. For instance, if an outer block is removed, the depth of the break instructions within nested blocks must be updated to reflect the new enclosing block depth. This is a significant challenge for rewriting tools, as it requires the analysis of the control-flow graph to determine the enclosing block depth for each break instruction.

Notice that, WebAssembly’s control flow design adheres to a Control Flow Integrity (CFI) policy. CFI is a security mechanism that limits a program’s control flow to a specified set of valid targets, thereby preventing arbitrary jumps [48]. Thus, even when a WebAssembly program originates from potentially untrustworthy sources, CFI policy theoretically guarantees the prevention of arbitrary jumps to random code locations.

2.1.5 Security and Reliability for WebAssembly

The WebAssembly ecosystem’s expansion needs robust tools to ensure its security and reliability. Numerous tools, employing various strategies to detect vulnerabilities in WebAssembly programs, have been created to meet this need. This section presents a review of the most relevant tools in this field, focusing on those capable of providing security guarantees for WebAssembly binaries.

Static analysis: SecWasm[49] uses information control-flow checking to identify vulnerabilities in WebAssembly binaries. Similarly, Wasmati[50] employs code property graphs for this purpose. Wasp[51] leverages concolic execution to identify potential vulnerabilities in WebAssembly binaries. CT-Wasm[52], verifies the constant time implementation of cryptographic algorithms in WebAssembly. Similarly, Vivienne applies relational Symbolic Execution (SE) to WebAssembly binaries in order to reveal vulnerabilities in cryptographic implementations[53]. While these tools emphasize specific strategies, others adopt a more holistic approach. For example, both Wassail[54] and WasmA[55] provide a comprehensive static analysis framework for WebAssembly binaries. Yet, static analysis tools may have limitations. For instance, a newly, semantically equivalent WebAssembly binary may be generated from the same source code bypassing or breaking the static analysis [56]. Thus, there may be a lack of subjects to evaluate the effectiveness of these tools.

Dynamic analysis: Dynamic analysis involves tools such as TaintAssembly[57], which conducts taint analysis on WebAssembly binaries. Furthermore, Stiévenart and colleagues have developed a dynamic approach to slicing WebAssembly programs based on Observational-Based Slicing (ORBS)[58, 59]. This technique aids in debugging, understanding programs, and conducting security analysis. However, Wasabi[60] remains the only general-purpose dynamic analysis tool for WebAssembly binaries, primarily used for profiling, instrumenting, and debugging

WebAssembly code. Similar to static analysis, these tools typically analyze software behavior during execution, making them inherently reactive. In other words, they can only identify vulnerabilities or performance issues while pseudo-executing input WebAssembly programs. Thus, facing an important limitation on overhead for real-world scenarios.

Protecting WebAssembly binaries and runtimes: The techniques discussed previously are primarily focused on reactive analysis of WebAssembly binaries. However, there exist approaches to harden WebAssembly binaries, enhancing their secure execution, and hardening the security of the entire execution runtimes ecosystem. For instance, Swivel[61] proposes a compiler-based strategy designed to eliminate speculative attacks on WebAssembly binaries in Function-as-a-Service (FaaS) platforms. Similarly, Kolosick and colleagues [62] modify the Lucet compiler to use zero-cost transitions, eliminating the performance overhead of SFI guarantees implementation. In addition, WaVe[63] introduces a mechanized engine for WebAssembly that facilitates differential testing. WaVe can be employed to detect anomalies in engines running Wasm-WASI programs. Much like static and dynamic analysis tools, these tools may suffer from a lack of WebAssembly inputs, which could affect the measurement of their effectiveness.

WebAssembly malware: Since the introduction of WebAssembly, the Web has consistently experienced an increase in cryptomalware. This rise primarily stems from the shift of mining algorithms from CPUs to WebAssembly, a transition driven by notable performance benefits [64]. Tools such as MineSweeper[23], MinerRay[24], and MINOS[25] employ static analysis with machine learning techniques to detect browser-based cryptomalware. In addition, SEISMIC[26], RAPID[27], and OUTGuard[28] leverage dynamic analysis techniques to achieve a similar objective. Moreover, VirusTotal¹⁵, a tool incorporating over 60 commercial antivirus systems as black-boxes, is capable of detecting cryptomalware in WebAssembly binaries. However, obfuscation studies have exposed their shortcomings, revealing an almost unexplored area for WebAssembly that threatens malware detection accuracy. In concrete, Bhansali et al. seminal work[65] demonstrate that cryptomining algorithm's source code can evade previous techniques through the use of obfuscation techniques.

2.1.6 Open challenges

Despite progress in WebAssembly analysis, numerous challenges remain. WebAssembly, though deterministic and well-typed by design, is susceptible to a variety of security threats. First, most existing WebAssembly research is reactive, focusing on detecting and fixing vulnerabilities already reported. This approach

¹⁵<https://www.virustotal.com>

leaves WebAssembly binaries and runtime implementations potentially open to unidentified attacks. Second, side-channel attacks present a significant risk. Genkin et al., for example, illustrated how WebAssembly could be manipulated to extract data via cache timing-side channels [19]. Furthermore, research conducted by Maisuradze and Rossow demonstrated the potential for speculative execution attacks on WebAssembly binaries [20]. Rokicki et al. disclosed the possibility for port contention side-channel attacks on WebAssembly binaries in browsers [15]. Finally, the binaries themselves may be inherently vulnerable. For example, studies by Lehmann et al. and Stiévenart et al. suggested that flaws in C/C++ source code could infiltrate WebAssembly binaries [17, 18].

2.2 Software diversification

Software diversification involves the synthesis, reuse, distribution, and execution of different, functionally equivalent programs so-called software variants. As outlined in Baudry et al.’s survey [66], Software Diversification falls into five usage categories: reusability [67], performance [68], fault tolerance [69], software testing [70], and security [30]. Our work specifically contributes to the last two categories. Based on the works of Cohen et al. [30], Forrest et al. [31], Jackson et al. [71] and Baudry et al. [66], this section presents core concepts and related works.

Software variants refer to functionally equivalent versions of an original program, produced through Software Diversification at various stages of the software lifecycle, from dependencies (coarse-grained) to machine code levels (fine-grained). The main goal of Software Diversification is to increase the cost of exploitation by making software less predictable. Diversification may be natural [66] or automatic [72]. Natural diversity refers to the side effect of humans creating software variants using different programming languages, compilers, and operating systems [66], all of which adhere to the initial requirements. The software market and competition typically address the creation of natural diversity. For example, Firefox and Chrome web browsers demonstrate natural diversity due to their practical differences in implementation and performance, despite serving the same purpose. This logic extends to operating systems, database engines, virtual machines, and application servers [66]. Natural diversity significantly aids in system security, as different variants are not susceptible to the same vulnerabilities [73, 74]. Unlike N-Version programming [75], natural diversity organically emerges over decades. In other words, while it does not require the allocation of additional human efforts, natural diversity cannot be automatically generated. This is because it is a side effect of the software development process. Given that WebAssembly is a relatively new technology, natural diversity is presently not a feasible option. Hence, for WebAssembly, feasible options are systematic and automatic diversification approaches.

2.2.1 Automatic generation of software variants

The concept of automatic software variants starts with Randell’s 1975 work [76], which put forth the notion of artificial fault-tolerant instruction blocks. Artificial Software Diversification, as proposed by Cohen and Forrest in the 1990s [30, 31], gets its development through rewriting strategies. These strategies consist of rule sets for modifying software components to create functionally equivalent, yet distinct, programs. Rewriting strategies typically take the form of tuples: $\text{instr1} \Rightarrow (\text{instr2}, \text{instr3}, \dots)$, where instr represents the original code and $(\text{instr2}, \text{instr3}, \dots)$ denotes the functionally equivalent code.

Rewriting strategy: The automatic creation of Software Diversification commences with rewriting rules. A rewriting rule refers to a functionally equivalent substitution for a code segment, manually written. These rules can be applied at varying levels, from coarse to fine-grained. This can range from the program dependencies level [77] to the instruction level [72]. For example, Cleempot et al. [78] and Homescu et al. [79] inject NOP instructions to yield statically varied versions at the instruction level. Here, the rewriting rule is represented as $\text{instr} \Rightarrow (\text{nop instr})$, signifying a `nop` operation preceding the instruction.

Instruction Reordering: This strategy reorders instructions in a program. For example, variable declarations may change if compilers reorder them in the symbol tables. This prevents static examination and analysis of parameters and alters memory locations. In this area, Bhatkar et al. [80, 81] proposed the random permutation of variable and routine order for ELF binaries. Such strategies are not implemented for WebAssembly to the best of our knowledge.

Adding, Changing, Removing Jumps and Calls: This strategy generates program variants by adding, changing, or removing jumps and calls in the original program. Cohen [30] primarily illustrated this concept by inserting random jumps in programs. Pettis and Hansen [82] suggested splitting basic blocks and functions for the PA-RISC architecture, inserting jumps between splits. Similarly, Crane et al. [83] de-inlined basic blocks of code as an LLVM pass. In their approach, each de-inlined code transforms into semantically equivalent functions that are randomly selected at runtime to replace the original code calculation. On the same topic, Bhatkar et al. [81] extended their previous approach [80], replacing function calls with indirect pointer calls in C source code, allowing post-binary reordering of function calls. In the WebAssembly context, the most analogous work is wobfuscator [84]. Wobfuscator, a JavaScript obfuscator, substitutes JavaScript code with WebAssembly code, e.g., numeric calcula. This strategy effectively uses the interleaving of calls between JavaScript and WebAssembly to provide JavaScript variants.

Program Memory and Stack Randomization: This strategy alters the layout of programs in the host memory. Additionally, it can randomize how a

program variant operates its memory. The work of Bhatkar et al. [80, 81] proposes to randomize the base addresses of applications and library memory regions in ELF binaries. Tadesse Aga and Autin [85], and Lee et al. [86] propose a technique to randomize the local stack organization for function calls using a custom LLVM compiler. Younan et al. [87] suggest separating a conventional stack into multiple stacks where each stack contains a particular class of data. On the same topic, Xu et al. [88] transforms programs to reduce memory exposure time, improving the time needed for frequent memory address randomization. This makes it very challenging for an attacker to ignore the key to inject executable code. This strategy disrupts the predictability of program execution and mitigates certain exploits such as speculative execution. No work has been found that explicitly applies this strategy to WebAssembly.

ISA Randomization and Simulation: This strategy involves using a key to cypher the original program binary into another encoded binary. Once encoded, the program can only be decoded at the target client, or it can be interpreted in the encoded form using a custom virtual machine implementation. This technique is strong against attacks involving code inspection. Kc et al. [89], and Barrantes et al. [90] proposed seminal works on instruction-set randomization to create a unique mapping between artificial CPU instructions and real ones. On the same topic, Chew and Song [91] target operating system randomization. They randomize the interface between the operating system and the user applications. Couroussé et al. [92] implement an assembly-like DSL to generate equivalent code at runtime in order to increase protection against side-channel attacks. Their technique generates a different program during execution using an interpreter for their DSL. Generally, *ISA randomization and simulation* usually faces a performance penalty, especially for WebAssembly, due to the decoding process as shown in WASMixer evaluation [56].

Code obfuscation: Code obfuscation can be seen as a simplification of *ISA randomization*. The main difference between encoding and obfuscating code is that the former requires the final target to know the encoding key while the latter executes as is in any client. Yet, both strategies aim to tackle static reverse engineering of programs. In the context of WebAssembly, Romano et al. [84] proposed an obfuscation technique, wobfuscator, for JavaScript in which part of the code is replaced by calls to complementary WebAssembly functions. Yet, wobfuscator targets JavaScript code, not WebAssembly binaries.

Enumerative synthesis: Enumerative synthesis is a fully automated and systematic approach to generate program variants. It examines all possible programs specific to a given language. The process of enumerative synthesis commences with a piece of input program, typically a basic block. Incrementally, using a defined grammar, it generates all programs of size n . A generated program is then checked for equivalence to the original program, either by using a test suite or a theorem solver. If the generated variant is proved, it

is added to the variants collection. The procedure continues until all potential programs have been explored. This approach proves especially effective when the solution space is relatively small or can be navigated efficiently. Jacob and colleagues [93] implemented this strategy for x86 programs. They named this technique superdiversification, drawing parallels to superoptimization [94]. Since this strategy fully explores a program’s solution space, it contains the aforementioned strategies as special cases. The application of enumerative synthesis to WebAssembly has not been explored.

2.2.2 Equivalence Checking

Equivalence checking between program variants is a vital component for any program transformation task, ranging from checking compiler optimizations [95] to the artificial synthesis of programs discussed in this chapter. It proves that two pieces of code or programs are functionally equivalent [96]. We can roughly simplify the checking process with the following property: two programs are deemed equivalent if they generate identical outputs when given identical inputs from a closed collection of inputs [97]. We adopt this definition of *functional equivalence modulo input* throughout this dissertation. In Software Diversification, equivalence checking seeks to preserve the original functionality of programs while varying observable behaviors. Two programs, for instance, can differ statically and still compute the same result. We outline two methods to check variant equivalence: by construction and prove-driven equivalence checking.

Equivalence checking by construction: The equivalence property is often guaranteed by construction. Cleemput et al. [78] and Homescu et al. [79], for example, design their transformation strategies to generate semantically equivalent program variants. However, developer errors can occur in this process, necessitating further validation. The test suite of the original program can serve as a check for the variant. If the program variant passes the test suite [98], it can be considered equivalent to the original. However, this technique is limited by the need for a preexisting test suite and does not give guarantees. An alternative method for checking program equivalence involves the use of fuzzers [99]. Fuzzers randomly generate inputs that yield different observable behaviors. If two inputs produce a different output in the variant, the variant and the original program are not equivalent. The primary limitation for fuzzers is that the process is notably time-consuming and necessitates manual introduction of oracles. Recent advances in machine learning have prompted researchers to investigate the use of neural networks for verifying program equivalence. One such example is the work of Zhang and colleagues [100], which produces reference oracles and test cases by using Large Language Models. Although this method appears effective, it only achieves an accuracy rate of 88%, which is insufficient for total verification.

Formal checking: In the absence of a test suite or a technique that inherently

implements the equivalence property, the works mentioned earlier use theorem solvers (SMT solvers) [101] to prove equivalence of program variants. The central idea for SMT solvers is to convert the two code variants into mathematical formulas. The SMT solver then checks for counter-examples [102]. When it finds a counter-example, there is an input for which the two mathematical formulas yield different outputs. The primary limitation of this technique is that not all algorithms can be translated into a mathematical formula, such as loops. Nevertheless, this technique is frequently used for checking no-branching programs like basic block and peephole replacements [103].

2.2.3 Variants deployment

Program variants, once generated and verified, may be utilized in two primary scenarios: Randomization or Multivariant Execution (MVE) [71].

Randomization: In the context of our work, the term *Randomization* denotes a program’s ability to present different variants to different clients. In this setup, a program, chosen from a collection of variants (referred to as the program’s variant pool), is assigned to a random client during each deployment. Jackson et al. [71] define the variant pool in Randomization as herd immunity, as vulnerable binaries can only affect a segment of the client community. El-Khalil and colleagues [104] suggest employing a custom compiler to generate varying binaries from the compilation process. They adapt a version of GCC 4.1 to partition a conventional stack into several component parts, termed multistacks. Similarly, Singhal and colleagues, propose Cornucopia [105]. Cornucopia generates multiple variants of a program by using different compiler flag combinations. Aga and colleagues [85], contributing to this discussion, propose the generation of program variants through the randomization of its data layout in memory. This method allows each variant to operate on the same data in memory but at different memory offsets. Randomization can also be applied to virtual machines and operating systems. On this note, Kc et al. [89] establish a unique mapping between artificial CPU instructions and actual ones, enabling the assignment of various variants to specific target clients. In a similar vein, Xu et al. [88] recompile the Linux Kernel to minimize the exposure time of persistent memory objects, thereby increasing the frequency of address randomization.

Multivariant Execution (MVE): Multiple program variants are composed into a single binary, known as a multivariant binary [106]. Each multivariant binary is randomly deployed to a client. Then, the multivariant binary executes its embedded program variants at runtime. These embedded variants can either execute in parallel to check for inconsistencies, or as a single program to randomize execution paths [80]. Bruschi and colleagues extend the concept of executing two variants in parallel, introducing non-overlapping and randomized memory layouts [107]. At the same time, Salamat et al. modifies a standard library

to generate 32-bit Intel variants. These variants have a stack that grows in the opposite direction, allowing for the detection of memory inconsistencies [108]. Davi and colleagues propose Isomeron, an approach for execution-path randomization [109]. Isomeron operates by simultaneously loading the original program and a variant. It then uses a coin flip to determine which copy of the program to execute next at the function call level. Previous works have highlighted the benefits of limiting execution to only two variants in a multivariant environment. Agosta and colleagues, as well as Crane and colleagues, used more than two generated programs in the multivariant composition, thereby randomizing software control flow at runtime [110, 83]. Both strategies have proven effective in enhancing security by addressing known vulnerabilities, such as Just-In-Time Return-Oriented Programming (JIT-ROP) attacks [111] and power side-channel attacks [112]. Lastly, only Voulimeneas et al. [113] have recently proposed a multivariant execution system that enhances security by parallelizing the execution of variants across different machines.

2.2.4 Measuring Software Diversification

Measuring Software Diversification presents a significant challenge. The size of the variant space does not necessarily correlate with a variant’s capacity to fulfill an objective such as hardening attacks [30]. Ideally, real scenarios would provide the most accurate measurement of diversification, e.g., demonstrating a variant’s effectiveness under specific attacks. However, such an approach is not always feasible, e.g., as previously discussed, Software Diversification is a preemptive strategy. Hence, a combination of static and dynamic metrics is required for measuring Software Diversification.

Static comparison of variants: Static metrics are used to measure diversity of programs without needing execution. The fundamental concept entails comparing variant source codes or binary codes to determine how diverse they are. Usually, comparing variants means defining a distance metric between programs. At the low-level of bytecode instructions, for example, these metrics include Levenshtein distance [114], and global alignments [33]. On the other hand, at the high-level of source code, these metrics often rely on Abstract Syntax Tree (AST) diffing, such as GUMtree-based distances [115]. Bostani et al. [116] illustrate the use of static distances in guiding the generation process of variants. They categorize the space of Android applications into malware and goodware. Then, they create malware variants by employing a static distance metric to approach the goodware group as closely as possible, thus successfully evading malware classifiers.

Dynamic comparison of variants: Static comparisons between variants inherently have limitations. For example, two variants may show differences at the source code level but exhibit identical behavior during execution. Take the addition of `nop` operations to a program as an instance. Despite source

code level differences, the variant and the original program execute identical instructions, leading to similar behaviors modulo input. Measuring Software Diversification primarily aims to demonstrate variant-specific observabilities. While static differences are observable, runtime information holds complementary relevance [117]. Therefore, dynamic metrics are essential to assess the diversity of variants. For instance, Forrest et al. [118] were pioneers in classifying program behaviors by analyzing their system call traces using n-grams profiling. Cabrera et al. utilized a global alignments approach to gauge the diversity of JavaScript bytecode traces within the Chrome browser [12]. Fang et al. proposed a method to counteract JavaScript obfuscation techniques used in malicious code, by analyzing dynamic information captured from V8 bytecode traces [119]. Dynamic metrics are primarily employed to cluster similar behaviors. Following the same logic, the diversity is greater when the difference between behaviors is larger. Notice that, dynamic metrics can be difficult due to the expense of program execution or the complication of required user interaction. On the other hand, malware programs, which usually do not require user interaction, are simpler to evaluate in controlled environments before actual deployment.

In the context of WebAssembly, there exist no explicit works on Software Diversification. Consequently, previous metrics have not been directly applied to measure diversification in WebAssembly binaries. However, in other domains, such as the analysis of WebAssembly binaries, several studies have employed static metrics. For example, VeriWasm quantifies attack-based patterns, stating that a WebAssembly binary is more secure with a lower pattern count [120]. This metric might potentially serve as a guide during variant generation. In the field of malware detection, MINOS [25] proposes transforming WebAssembly binaries into grayscale images. They then employ convolutional neural networks to identify malware, where an increased similarity to a malware image increases the probability of the binary being malware. Regarding the dynamic comparisons, Wang et al.'s study [26] profiles WebAssembly instructions during runtime to identify malicious behavior.

2.2.5 Offensive or Defensive assessment of diversification

Lundquist and colleagues [72] distinguish Software Diversification into two categories: Defensive and Offensive Diversification. On the one hand, Defensive Software Diversification introduces unpredictability in system behavior. By making software less predictable, defensive Software Diversification aims to proactively deter attacks, acting as a complementary strategy to other, more reactive, security measures. The majority of previously discussed works in this section contribute to defensive diversification. Yet, Software Diversification that aims to create diverse harmful programs is considered Offensive Diversification [121].

Offensive Diversification: Offensive Diversification is conceptually equal to Defensive Software Diversification. Yet, in an offensive context, one may apply diversification techniques to malware or other malicious codes to evade detection by security software [122]. For example, one might equate Offensive Diversification with Code obfuscation, if its purpose shifts from preventing reverse engineering by malicious actors, to evading detection by malware analysis systems.

Malicious actors may employ previously discussed diversification strategies to evade detection, including genetic programming [123]. Over time, these evasion techniques have evolved in both complexity and sophistication [124]. Chua et al. [125], for instance, suggested a framework for automatically obfuscating the source code of Android applications using method overloading, opaque predicates, try-catch, and switch statement obfuscation, resulting in multiple versions of identical malware. Moreover, machine learning approaches have been utilized to develop evasive malware [126], drawing on a corpus of pre-existing malware [116]. These methods aim to thwart static malware detectors, yet, more advanced techniques focus on evading dynamic detection mostly by employing throttling [127, 128].

The term Offensive Software Diversification may appear counterintuitive. Yet, such approaches measure the resilience and accuracy of security systems. This is an almost unexplored area in WebAssembly, posing a threat to malware detection accuracy. Specifically, only Bhansali et al.’s seminal work[65] has demonstrated that a cryptomining algorithm’s source code can evade pre-existing malware detection methods. More recently, Madvex [129] has sought to obfuscate WebAssembly binaries to achieve malware evasion, but this approach is limited to altering only the code section of WebAssembly binaries.

2.3 Open challenges for Software Diversification

As outlined in Section 2.1.6, our primary motivation for the contributions of this thesis is the open issues within the WebAssembly ecosystem. We see potential in employing Software Diversification to address them. Based on our previous discussion, we highlight several open challenges in the realm of Software Diversification for WebAssembly. First, WebAssembly, being an emerging technology, is still in the process of implementing defensive measures. The process of officially adopting a new defensive measure is inherently slow, making Software Diversification a potentially valuable strategy. Second, despite the abundance of related work on software diversity, its exploration in the context of WebAssembly remains limited. This thesis is first to investigate Software Diversity in depth for the emerging WebAssembly ecosystem. Third, both randomization and multivariant execution remain largely unexplored within the WebAssembly context. The deployment of Software Diversification in WebAssembly poses

unique challenges. WebAssembly ecosystems are remarkably dynamic. Web browsers and FaaS platforms serve as prime examples. In these environments, WebAssembly binaries are served millions of times simultaneously to the former, while new WebAssembly binaries are cold-spawned and executed upon each user request in the latter. Thus, designing practical Software Diversification for WebAssembly requires careful consideration of the deployment environment. Last but not least, research on malware detection, as discussed in Section 2.1.5, suggests that offensive diversification may assist in evaluating the resilience and accuracy of WebAssembly’s security systems.